

9200
RECEIVED
DEC 14 2004
Technology Center 2600

DT07 Rec'd PCT/PTO 08 JUL 2004

Tel: ++1-321-7284059
E-mail: msilaghi@cs.fit.edu
Mail: 1311 Harvard Circle #4,
Melbourne FL-32905, US

June 30, 2004

Application No. 09/647,300 Art Unit 2655
Mail Stop PCT,
Commissioner for Patents,
P.O.Box 1450, Alexandria, VA 22313-1450

This is a petition to your notice of abandonment of the Application No. 09/647,300, Art Unit 2655, that you mailed out on June 16, 2004

In your notice of abandonment it is argued that "No reply was received to your letter of 30 October 2003". I want to bring to your knowledge that I mailed to your patent office two answers to the corresponding notice, in December 2003 respectively January 2004, plus a faxed answer during January 2004 and several telephone inquiries about the status of my answers to supervisor Doris To.

I attach photocopies of the mailing receipts with the corresponding tracking numbers. I also attach the other documents that were asked in your letter, as I did in the previous answers:

- marked up version of the description with changes
- clean version of the description with changes
- fixed version of the claims
- summary
- disclosure of an article
- copy of the article

In the following I also include the content of my previous answer to you on January 17, that you might not have received:

Answer to your recent note concerning the patent application 09/647,300, Art Unit 2655, Examiner Daniel Demelash Abebe

This is an answer to your previous note related to the International application PCT/IB00/00189, US National stage 09/647,300, Art Unit 2655, Examiner Daniel Demelash Abebe.

- About the patent application RO-C99-00214 that the examiner cannot locate: according to my knowledge, a copy of it was transmitted to you

by the PCT office in Geneva (please verify the documents transmitted to you from PCT IB in Geneva).

- I also attach an Information disclosure statement filled with the changes that you recommended in your previous note. I actually attach several versions since I was not sure about different interpretations of your markings in the previous note.
- The description and claims were numbered and reformatted with double space.
- The two markups (on page 5 respectively 6) were reversed to the previous version. There was no markup on page 10. A substitute specification with amendments is submitted. The reference to the parent application is also added in the beginning of the specification. Both a clean and a marked-up version are included, 37 CFR 1.125(b)(2), 37 CFR 1.125(c).
- An amendment was made to each claim, by canceling and re-presenting/rewriting it. The formulation of the claims was done in the specified one sentence form. The translation is improved by correcting some mentioned spelling, grammar, and idiomatic errors.

Statement conforming 37 CFR 1.125(b)(1) *The substitute specification contains no new matter.*

I thank you in advance

Marius Călin Silaghi

Customer Copy
Label 11-B September 2002



ER 760593680 US



Post Office To Addressee

ORIGIN (POSTAL USE ONLY)			
PO ZIP Code 32905	Day of Delivery <input type="checkbox"/> Next <input checked="" type="checkbox"/> Second <input type="checkbox"/> Third	Flat Rate Envelope <input type="checkbox"/>	
Date In Mo. 12 Day 04 Year 2003	<input type="checkbox"/> 12 Noon <input checked="" type="checkbox"/> 3 PM	Postage \$ 1.30	
Time In <input type="checkbox"/> AM <input type="checkbox"/> PM	Military <input type="checkbox"/> 2nd Day <input type="checkbox"/> 3rd Day	Return Receipt Fee \$	
Weight lbs. 11 oz. 12	Int'l Alpha Country Code	COD Fee \$	Insurance Fee \$
No Delivery <input type="checkbox"/> Weekend <input type="checkbox"/> Holiday	Acceptance Clerk Initials	Total Postage & Fees \$	

DELIVERY (POSTAL USE ONLY)

Delivery Attempt	Time	Employee Signature
Mo. Day	<input type="checkbox"/> AM <input type="checkbox"/> PM	
Delivery Attempt	Time	Employee Signature
Mo. Day	<input type="checkbox"/> AM <input type="checkbox"/> PM	
Delivery Date	Time	Employee Signature
Mo. Day	<input type="checkbox"/> AM <input type="checkbox"/> PM	

CUSTOMER USE ONLY

PAYMENT BY ACCOUNT
Express Mail Corporate Acct. No.
Federal Agency Acct. No. or
Postal Service Acct. No.

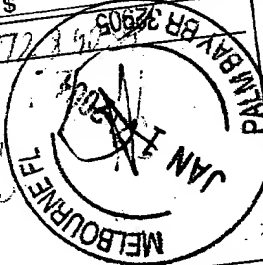
☐ WAIVER OF SIGNATURE (Domestic Only)
Additional merchandise insurance is valid if waiver of signature is requested.
I wish delivery to be made without obtaining signature of addressee or addressee's agent (if delivery employee judges that article can be left in secure location) and I authorize that delivery employee's signature constitutes valid proof of delivery.

NO DELIVERY
☐ Weekend ☐ Holiday

Customer Signature

FROM: (PLEASE PRINT)

MARY C SLOAN
PALM BAY CRIDE #4
MELBOURNE FL 32905



TO: (PLEASE PRINT)

MAIL STOP PCT, COMMISSIONER for Paleds
P.O. BOX 1450, ALEXANDRIA
VA 22313-1450

FOR PICKUP OR TRACKING CALL 1-800-222-1811

www.usps.com

PRESS HARD. You are making 3 copies.

U.S. Postal Service™ Delivery Confirmation™ Receipt

Postage and Delivery Confirmation fees must be paid before mailing.
Article Sent To: (to be completed by mailer)
MAIL STOP PCT, COMMISSIONER for Paleds
P.O. BOX 1450, ALEXANDRIA
VA 22313-1450

DELIVERY CONFIRMATION NUMBER: 0962 0960 0960 0000 0000 0000 0000 0000



POSTAL CUSTOMER:
Keep this receipt. For inquiries:
Access internet web site at
www.usps.com
or call 1-800-222-1811

CHECK ONE (POSTAL USE ONLY)

<input checked="" type="checkbox"/>	Priority Mail™ Service
<input type="checkbox"/>	First-Class Mail® parcel
<input type="checkbox"/>	Package Services parcel

(See Reverse)

PS Form 152, May 2002

RECEIVED

DEC 14 2004

Technology Center 2600

MARKUP, JUNE 129/2004

SPECIFICATION

3 TITLE OF THE INVENTION

Speech Recognition and Signal Analysis by Exact Fast Search of Subsequences with Maximal Confidence Measure

6

REFERENCE TO APPENDIX SUBMITTED ON CD

Not Applicable

9

CROSS-REFERENCE TO RELATED APPLICATION

This patent application has as parent application the patent application C99-00214/25.02.1999 registered with the State Office for Inventions and Trademarks (OSIM) in Bucharest, Romania. The present application is the US national stage of the international application PCT/IB00/00189 registered with the International Patent Office in Geneva.

15

BACKGROUND OF THE INVENTION

18 FIELD OF THE INVENTION

The invention relates to a common component of Speech Recognition, more particularly to the fields of Keyword Spotting and decoding, Segments Alignment for DNA and proteins, and Recognition of Objects in Images.

DESCRIPTION OF THE RELATED ART

This invention addresses the problem of keyword spotting (KWS) in unconstrained speech without explicit modeling of non-keyword segments (typically done by using filler HMM models or an ergodic HMM composed of context dependent or independent phone models without lexical constraints). Several methods (sometimes referred to as sliding model methods) tackling this type of problem have already been proposed in the past. E.g., they use Dynamic Time Warping (DTW) or Viterbi matching allowing relaxation of the (begin and endpoint) constraints. These are known to require the use of an appropriate normalization of the matching scores since segments of different lengths have then to be compared. However, given this normalization and the relaxation of begin/endpoints, straightforward Dynamic Programming (DP) is no longer optimal (or, in other words, the DP optimality principle is no longer valid) and has to be adapted, involving more memory and CPU. Indeed, at any possible ending time e , the match score of the best warp and start time b of the reference has to be computed (for all possible start times b associated with unpruned paths). Finally, this adapted DP quickly becomes even more complex (or intractable) for more advanced scoring criteria (such as the confidence measures mentioned below).

Work in the field of confidence level, and in the framework of hybrid HMM/ANN systems has shown that the use of accumulated local posterior probabilities (as obtained at the output of a multilayer perceptron) normalized by the length of the word segment (or, better, involving a double normalization over the number of phones and the number of acoustic frames in each phone) was yielding good confidence measures and good scores for the re-estimation of N -best hypotheses. However, so far the evaluation of such confidence measures

involved the estimation and rescoreing of N-best hypotheses.

KWS methods without filler models have in common the selection of a subsequence of the utterance to match the interesting keyword models. Let $X = \{x_1, x_2, \dots, x_n, \dots, x_N\}$ denote the sequence of acoustic vectors in which we want to detect a keyword, and let M be the HMM model of a keyword M and consisting of L states $\mathcal{Q} = \{q_1, q_2, \dots, q_\ell, \dots, q_L\}$. Assuming that M is matched to a subsequence $X_b^e = \{x_b, \dots, x_e\}$ ($1 \leq b \leq e \leq N$) of X , and that we have an implicit (not modeled) garbage/filler state q_G preceding and following M , one can define (approximate) the log posterior of a model M given a subsequence X_b^e as the average posterior probability along the optimal path, i.e.:

$$\begin{aligned}
-\log P(M|X_b^e) &\simeq \frac{1}{e-b+1} \min_{\forall Q \in M} -\log P(Q|X_b^e) \\
&\simeq \frac{1}{e-b+1} \min_{\forall Q \in M} \{-\log P(q^b|q_G) \\
&\quad - \sum_{n=b}^{e-1} [\log P(q^n|x_n) + \log P(q^{n+1}|q^n)] \\
&\quad - \log P(q^e|x_e) - \log P(q_G|q^e)\}
\end{aligned} \tag{1}$$

where $Q = \{q^b, q^{b+1}, \dots, q^e\}$ represents one of the possible paths of length $(e-b+1)$ in M , and q^n the HMM state visited at time n along Q , with $q^n \in \mathcal{Q}$. In this expression, q_G represents the garbage (filler) state which is simply used here as the non-emitting initial and final state of M . Transition probabilities $P(q^b|q_G)$ and $P(q_G|q^e)$ can be interpreted as the keyword entrance and exit penalties, but can be simply set to 1. Local posteriors $P(q_\ell|x_n)$ can be estimated using any of the known techniques: multi-gaussians, code-books, or as output values of a multilayer perceptron (MLP) used in hybrid HMM/ANN systems. For a specific sub-sequence X_b^e , expression (1) can easily be estimated by dynamic programming since the sub-sequence and the associated normalizing factor $(e-b+1)$ are given. However, in the

case of keyword spotting, this expression should be estimated for all possible begin/endpoint pairs $\{b, e\}$ (as well as for all possible word models), and we define the matching score of X

3 on M as:

$$S(M|X) = -\log P(M|X_{b^*}^{e^*}) \quad (2)$$

where the optimal begin/endpoints $\{b^*, e^*\}$, and the associated optimal path Q^* , are the

6 ones yielding the lowest average local posterior:

$$\langle Q^*, b^*, e^* \rangle = \operatorname{argmin}_{\{Q, b, e\}} \frac{-1}{e - b + 1} \log P(Q|X_b^e) \quad (3)$$

Of course, in the case of several keywords, all possible models will have to be evaluated.

9 A double averaging involving the number of frames per phone and the number of phones usually yields slightly better performance when used to rescore N-best candidates:

$$\langle Q^*, b^*, e^* \rangle = \quad (4)$$

$$\operatorname{argmin}_{\{Q, b, e\}} \frac{-1}{J} \sum_{j=1}^J \left(\frac{1}{e_j - b_j + 1} \sum_{n=b_j}^{e_j} \log P(q_j^n | x_n) \right)$$

where J represents the number of phones in the hypothesized keyword model and q_j^n the

hypothesized phone q_j for input frame x_n . However, given the time normalization and

15 the relaxation of begin/endpoints, straightforward DP is no longer optimal and has to be adapted, usually involving more memory and CPU.

Filler-based KWS need a simpler decoding step. Although various solutions have been

18 proposed towards the direct optimization of (2), most of the keyword spotting approaches

today prefer to preserve the optimality and simplicity of Viterbi DP by modeling the complete

input and explicitly or implicitly modeling non-keyword segments by using so called filler or

21 garbage models as additional reference models. In this case, we assume that non-keyword

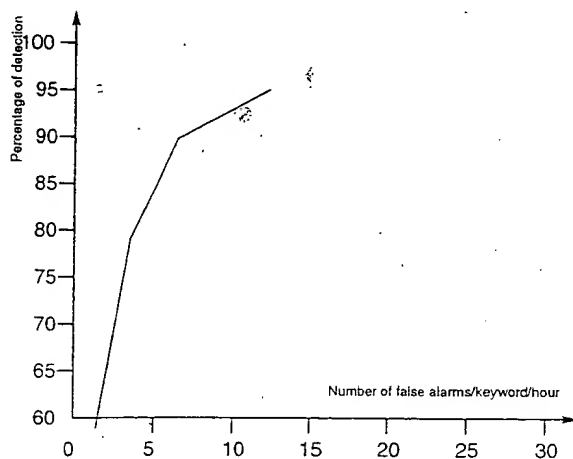


Figure 2: ROC using criterion (4) (double normalization), on 242 BREF test sentences containing 100 keywords selected at random.

cific tuning) appear to be particularly competitive to other alternative approaches.

7. ACKNOWLEDGMENTS

We thank Giulia Bernardis for helping us with the BREF database and providing us with a neural network trained for this task. We also acknowledge the useful discussions with Prof. Boi Faltings.

8. REFERENCES

- [1] Bernardis, G. and Boulard, H., "Improving posterior-based confidence measures in hybrid HMM/ANN speech recognition systems," *Proceedings of Intl. Conf. on Spoken Language Processing* (Sydney, Australia), pp. 775-778, 1998.
- [2] Boulard, H. and Morgan, N., *Connectionist Speech Recognition - A Hybrid Approach*, Kluwer Academic Publishers, 1994.
- [3] Boulard, H., D'hoore, B., and Boite, J.-M., "Optimizing recognition and rejection performance in wordspotting systems," *Proc. of IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing* (Adelaide, Australia), pp. I:373-376, 1994.
- [4] Bridle, J.S., "An efficient elastic-template method for detecting given words in running speech," *Proc. of the Brit. Acoust. Soc. Meeting*, pp. 1-4, April 1973.
- [5] Lamel, L.F., Gauvain, J.-L., and Eskénazi, M., "BREF, a large vocabulary spoken corpus for French," *Proceedings of Eurospeech'91*, pp. 505-508, 1991.
- [6] Rohlicek, J.R., "Word spotting," in *Modern Methods of Speech Processing*, R.P. Ramachandran and R. Mammone (Eds.), Kluwer Academics Publishers, pp. 123-157, 1995.
- [7] Rose, R.C. and Paul, D.B., "A hidden Markov model based keyword recognition system," *Proc. of ICASSP'90*, pp. 129-132, 1990.
- [8] Silaghi, M.-C. and Boulard H., "Posterior-Based Keyword Spotting Approaches Without Filler Models," *Swiss Federal Institute of Technology Lausanne (EPFL)*, Technical Report, 1999.
- [9] Sukkar, R.A. and Lee, C.-H., "Vocabulary independent discriminative utterance verification for nonkeyword rejection in subword based speech recognition," *IEEE Trans. on Speech and Audio Processing*, vol. 4, no. 6, pp. 420-429, 1996.
- [10] Williams, G. and Renals, S., "Confidence measures for hybrid HMM/ANN speech recognition," *Proceedings of Eurospeech'97*, pp. 1955-1958, 1997.
- [11] Wilpon, J.G., Rabiner, L.R., Lee C.-H., and Goldman, E.R., "Application of hidden Markov models of keywords in unconstrained speech," *Proc. of ICASSP'89*, pp. 254-257, 1989.

segments are modeled by extraneous garbage models/states q_G (and grammatical constraints ruling the possible keyword/non-keyword sequences).

3 [It is sufficient to consider only the case of detecting one keyword] *Let*
us consider only the case of detecting one keyword per utterance at a time. In this case,
the keyword spotting problem amounts at matching the whole sequence X of length N onto
6 an extended HMM model \overline{M} consisting of the states $\{q_G, q_1, \dots, q_L, q_G\}$, in which a path
(of length N) is denoted $\overline{Q} = \{\overbrace{q_G, \dots, q_G}^{b-1}, q^b, q^{b+1}, \dots, q^e, \overbrace{q_G, \dots, q_G}^{N-e}\}$ with $(b-1)$ garbage states
 q_G preceding q^b and $(N-e)$ states q_G following q^e , and respectively emitting the vector
9 sequences X_1^{b-1} and X_{e+1}^N associated with the non-keyword segments.

Given some estimation of $P(q_G|x_n)$ (e.g., using probability density functions trained on
non keyword utterances), the optimal path \overline{Q}^* (and, consequently b^* and e^*) is then given

12 by:

$$\begin{aligned} \overline{Q}^* &= \underset{\forall \overline{Q} \in \overline{M}}{\operatorname{argmin}} -\log P(\overline{Q}|X) \\ &= \underset{\forall \overline{Q} \in \overline{M}}{\operatorname{argmin}} \{-\log P(Q|X_b^e) \\ &\quad - \sum_{n=1}^{b-1} \log P(q_G|x_n) - \sum_{n=e+1}^N \log P(q_G|x_n)\} \end{aligned} \quad (5)$$

15 which can be solved by straightforward DP (since all paths have the same length). The main
problem of filler-based keyword spotting approaches is then to find ways to best estimate
18 $P(q_G|x_n)$ in order to minimize the error introduced by the approximations. Sometimes this
value was defined as the average of the N best local scores while, in other approaches, this
value is generated from explicit filler HMMs. However, these approaches will usually not
21 lead to the optimal solution given by (2).

BRIEF SUMMARY OF THE INVENTION

The invention belongs to the technical domain of decoding, classification, alignment and
3 matching of data.

The invention introduces a new method performing tasks in keyword spotting in utter-
ances, detection of subsequences in chains of organic matter (DNA and proteins) and recog-
6 nition of objects in images. The proposed methods search in an optimized way the matching
that maximizes, over all the possible matchings, certain confidence measures based on nor-
malized posteriors. Three such confidence measures are used, two existed in previous work
9 in Speech Recognition, and the third one is a new one.

Application fields for this invention are: man-machine interfaces (using speech recogni-
tion; ex: control systems, banking, flight services, etc), coordination systems (for industrial
12 robots and automata) and development systems for pharmaceutic products.

BRIEF DESCRIPTION OF THE SEVERAL VIEWS OF THE DRAWINGS

15 Not Applicable

DETAILED DESCRIPTION OF THE INVENTION

18 [The present invention introduces a fast iterative method,] = *In the following,*
we show that it is possible to define an iterative process, = referred to as Iterating Viterbi De-
coding (IVD) with good/fast convergence properties, estimating the value of $P(q_G|x_n)$ such
21 that straightforward DP (5) yields exactly the same segmentation (and recognition results)
than (3). While the same result could be achieved through a modified DP in which all pos-

sible combinations (all possible begin/endpoints) would be taken into account, the method proposed below is much more efficient (in terms of both CPU and memory requirements).

3 Compared to previously devised sliding model methods the first method proposed here is based on: (A) A matching score defined as the average observation probability (posterior) along the most likely state sequence. It is indeed believed that local posteriors are more
6 appropriate to the task. (B) The iteration of a Viterbi decoding [algorithm] \Rightarrow , which does not require scoring for all begin/endpoints or N-best rescoring, and which can be proved to (quickly) converge to the optimal (from the point of view of the chosen scoring functions)
9 solution without requiring any specific filler models, using straightforward Viterbi alignments (similar to regular filler-based KWS, but for some versions at the cost of a few iterations).

The IVD method is based on a similar criterion as the filler based approaches (5), but
12 rather than looking for explicit (and empirical) estimates of $P(q_G|x_n)$ we aim at mathematically estimating its value (which will be different and adapted to each utterance) such that solving (5) is equivalent to solving (3). Thus, we perform an iterative estimation of
15 $P(q_G|x_n)$, such that the segmentation resulting of (5) is the same than what would be obtained from (3). Defining $\varepsilon_t = -\log P(q_G|x_n)$ at iteration t , the proposed method can be summarized as follows:

- 18 1. Start the first iteration, $t = 0$, from an initial value $\varepsilon_0 = \Pi$ (it is actually proven that the iterative process presented here will always converge to the same solution, in more or less cycles with the worst case upper bound of N iterations, independently of this
21 initialization, e.g., with Π equal with a cheap estimation of the score of a match).

In one of the developed versions, ε_0 is initialized to $-\log$ of the maximum of the local

probabilities $P(q_k|x_n)$ for each frame x_n .

3 An alternative choice is to initialize ε_0 to a pre-defined threshold score, T , that expression (1) should reach to declare a keyword matching (see step 4 below). In this last case, if $\varepsilon_1 > \varepsilon_0$ at the first iteration, then we can (as proven) directly infer that the match will be rejected, otherwise it will be accepted.

6 2. Given the estimate ε_t of $P(q_G|x_n)$ at current iteration t , find the optimal path $\langle \overline{Q}_t, b_t, e_t \rangle$ according to (5) and matching the complete input.

3. Estimate the value of ε_{t+1} to be used in the next iteration as the average of the local
9 posteriors along the optimal path Q_t (matching the $X_{b_t}^{e_t}$ resulting of (5) on the keyword model) i.e.:

$$\varepsilon_{t+1} = -\frac{1}{(e_t - b_t + 1)} \log P(Q_t|X_{b_t}^{e_t}) \quad (6)$$

12 4. Increment t and return to (2) iterating until convergence is detected. If we are not interested in the optimal segmentation, this process could also be stopped as soon as it reaches a ε_{t+1} lower than a (pre-defined) minimum threshold, T , below which we can
15 declare that a keyword has been detected.

Correctness and convergence proof of this process and generalization to other criteria, are available: each IVD iteration (from the second iteration) will decrease the value of ε_t , and the
18 final path yields the same solution than (3). The above method has a very good experimental convergence speed (3-5 iterations in our tests). For one version of IVD (when ε_0 is initialized using the acceptance threshold, T), the detection is decided after one single step.

A version with the same effort but suboptimal results is proposed in the following paragraph. Let $T(\overline{M}, X)$ be a matrix holding the HMM emission probabilities for an utterance X whose time-frames define the columns, and where the states of the hypothesized word W define the rows. When using the standard DP, one computes for each element of the matrix $T(\overline{M}, X)$ at frame k of X and state s of \overline{M} three values: S_{ks} , L_{ks} and C_{ks} , where S_{ks} corresponds to the sum of the entries on the optimal path that leads to the entry, L_{ks} holds the length of the optimal path computed so far, and C_{ks} is the estimation of the cost on the optimal expanded path. By a path leading to an entry $T(k, s)$ we mean a sequence of entries in the table T , such that there is exactly an entry for each time frame $t \leq k$. At each entry $T(k, s)$, DP selects a locally optimal path noted P_{ks} . At each step k , we consider all pairs of entries of table $T(\overline{M}, X)$ of type $T(k, s)$, $T(k-1, t)$. We update for each such pair, the current cost C_{ks} (initially ∞), by comparing it with the alternative given by:

$$S_{ks} = S_{(k-1)t} - \log p(s|x_k)p(s|t)$$

$$L_{ks} = L_{(k-1)t} + 1, \forall t > 0, t \leq L$$

$$C_{ks} = \frac{S_k}{L_k} \tag{7}$$

wanting to have at step k the path P_{ks} from the paths $P_{(k-1)t}$ that minimizes C_{NL} . With DP, one will choose the P_{ks} with minimal C_{ks} .

This version can yield suboptimal results since the optimality principle is not respected by the expression 7. The optimality principle of Dynamic Programming requires that the path to the frame $k-1$ that minimizes C_{NL} , also minimizes C_{ks} for an entry at frame k of table $T(\overline{M}, X)$.

Another technique that is suboptimal in time and/or quality is obtained from the previous one adopting a beam-search approach and a set of safe prunings. The Dynamic Programming
3 can be viewed as a set of safe prunings that are applied at each entry of the DP table and has the property that only one alternative is maintained. Dynamic Programming cannot be used, since the principle of optimality is not respected. The following types of safe pruning
6 that can be done are introduced by the present invention. Within the current invention we found a set of safe prunings as follows: we have proved that if at a frame a we have two paths P'_a and P''_a with $S''_a < S'_a$ and $L'_a < L''_a$, then at no frame $c \geq a$ will a path P''_c be forsaken for
9 a path P'_c if $P'_a \subset P'_c$, $P''_a \subset P''_c$ and $P'_c \setminus P''_c \equiv P''_c \setminus P'_c$. We will note the order relation as $P''_a \prec P'_a$. We have further shown that a path P' may be safely discarded only when we know a lower cost one, P'' .

$$12 \quad P' \prec P'' \Rightarrow C'_k < C''_k \quad (8)$$

Thus, the method described in following method computes $S(M, X)$ and Q^* from equation (3). By ordering the set of paths, according to Equation 8, we only need to check the
15 step (1.1) of the following method up to the eventual insertion place. The last paths are candidates for pruning in step (1.2). In order for the pruning to be acceptable, we will prune only paths that were too long on the last state. An additional counter for each path is
18 needed for storing the state length. This counter is reset when an entry from another row is added and is incremented at each advance with a frame. The following steps detail this method for a model W and an utterance X :

- 21 a) Initialize all elements of a matrix, $\text{SetOfPaths}(1..N, 1..K)$, to \emptyset
- b) For all frames from 1 to N , for all states from 1 to K , for all candidates p_i in

SetOfPaths(frame-1, 1..K):

- For all p_j in SetOfPaths[frame, state], if $p_i \prec p_j$ then delete p_j (1.1), and if $p_j \prec p_i$ then continue step b) (1.2)
 - Insert p_i in SetOfPaths[frame, state]
- c) Select SetOfPaths[frame, K] as the best of the candidates

The next method builds on the previous technique and is a fast procedure for maximizing a more complex confidence measure that yields better results in practice. The corresponding confidence measure is defined as:

$$\frac{1}{NVP} \sum_{h_i \in VP} \frac{\sum_{pst \in h_i} -\log(pst)}{length(h_i)} \quad (9)$$

where NVP stands for the number of visited phonemes and VP stands for the set of visited phonemes. An average is computed over all posteriors pst of the emission probabilities for the time frames matched to the visited phoneme h_i . The function $length(h_i)$ gives the number of time frames matched against h_i . This method uses a breath first Beam Search algorithm. It exploits a set of reduction rules and certain normalizations. For the state q_G , in this method, the logarithm of the emission posterior is equal with zero. For each frame e and for each state s , the set of paths/probabilities of having the frame e in the state s is computed as the first \mathcal{N} maxima (\mathcal{N} can be finite) of the confidence measure for all paths in HMM \overline{M} of length e and ending in the state s . The paths that according to the reduction rules will loose the final race when compared with another already known path, will be deleted as well. Let us note a_1, p_1, l_1 , respectively a_2, p_2 and l_2 the confidence measure for the previously visited phonemes, the posterior in the current phoneme and the length in the

current phoneme for the path Q_1 , respectively the path Q_2 . The rules that can be used for the reduction of the search space by discarding a path Q_1 for a path Q_2 are in this case any

3 of the next ones:

1. $l_2 \geq l_1$, $A > 0$, $B \leq 0$ and $L_c^2 A + L_c B + C \geq 0$
2. $l_2 \geq l_1$, $A \geq 0$, $B \geq 0$ and $C \geq 0$
- 6 3. $l_2 \geq l_1$, $A \leq 0$, $C \geq 0$ and $L^2 A + LB + C \geq 0$
4. $l_2 \geq l_1$, $A = 0$, $B < 0$ and $LB + C \geq 0$

where $A = a_1 - a_2$, $B = (a_1 - a_2)(l_1 + l_2) + p_1 - p_2$, $C = (a_1 - a_2)l_1 l_2 + p_1 l_2 - p_2 l_1$, $L =$

9 $L_{max} - \max\{l_1, l_2\}$, $L_c = -B/2A \geq 0$ and L_{max} is the maximum acceptable length for a

phoneme. By discarding paths only if one of the above rules is satisfied, the optimum defined by the confidence measure with double normalization can be guaranteed, if no phone may be

12 avoided by the HMM M . Any HMM may be decomposed in HMMs with this quality. The

4-th rule is included in the 3-rd and its test is useless if the last one was already checked.

The first test, $l_2 \geq l_1$ tells us if Q_2 has chances to eliminate Q_1 , otherwise we will check

15 if Q_1 eliminates Q_2 . These tests were inferred from the conditions of maintaining the final

maximal confidence measure while reduction takes place. In order to use the method of double normalization without decomposing HMMs that skip some phonemes, the previous

18 rules are modified taking into account the number of visited phonemes for any path F_1

respectively F_2 and the number of phonemes that may follow the current state. A simplified

test can be:

- 21 • $l_2 \geq l_1$, $A \geq 0$, $p_1 \geq p_2$ respectively $F_2 \geq F_1$ for the HMMs that skips phonemes.

This test is weaker than the 2nd reduction rule. For example a path is eliminated by a second path if the first one has an inferior confidence measure (higher in value) for the the previous phonemes, a shorter length and the minus of the logarithm of the cumulated posterior in the current phoneme also inferior (higher in value) to that of the second one. An additional confidence measure based on the maximal length, L_{max} , and on the maximum of the minus of the logarithm of the cumulated and normalized posterior in phoneme, P_{max} , can be used in order to limit the number of stored paths.

- $p > L_{max}P_{max}$ in any state
- $\frac{p}{l} > P_{max}$ at the output from a phoneme

where p and l are the values in the current phoneme for the minus of the logarithm of cumulated posterior and for the length of the path that is discarded. These tests allow for the elimination of the paths that are too long without being outstanding, respectively of the paths with phonemes having unacceptable scores, otherwise compensated by very good scores in other phonemes. If \mathcal{N} is chosen equal with one, the aforementioned rules are no longer needed, but always we propagate the path with the maximal current estimation of the confidence measure. The obtained results are very good, even if the defined optimum is guaranteed for this method only when \mathcal{N} is bigger than the length of the sequence allowed by L_{max} or of the tested sequence. The same approach is valid for the simple normalization, where the HMM for the searched word will be grouped into a single phoneme.

The present invention can exploit a newly designed a confidence measure, version named Real Fitting, that represents differently the exigencies of the recognition. Since the phonemes and the absent states can be modeled by the used HMMs, we find it interesting to request the

fitting of each phoneme in the model with a section of the sequence. Therefore, we measure the confidence level of a subsequence as being equal with the maximum over all phonemes of the minus of the logarithm of the cumulated posterior of the phone, normalized with its length:

$$\max_{phonem \in Visited\ Phonems} \frac{\sum_{phonem} -\log(posterior)}{phonem\ length} \quad (10)$$

The rule that may be used in this framework for the reduction of the number of visited paths is:

- Q_2 is discarded in favor of another path Q_1 if the confidence measure of the Real Fitting for the previous phonemes is inferior (higher in value) for Q_2 compared with Q_1 , and if $p_1 \leq p_2$ and $l_2 \leq l_1$.

where p_1 , l_1 , respectively p_2 , l_2 represent the minus of the logarithm of the cumulated posterior respectively the number of frames in the current phoneme for the path Q_1 respectively Q_2 . Similarly to the previous method, the set of visited paths can be pruned by discarding those where:

- $p > L_{max}P_{max}$ in any state
- $\frac{p}{l} > P_{max}$ at the output from a phoneme

where p and l are the values in the current phoneme for the minus of the logarithm of the cumulated posterior and for the length of the path that is discarded. We recall that the meaning of the constants are the maximal length L_{max} , respectively the accepted maxima of the minus of the logarithm of the cumulated and normalized posterior in phoneme, P_{max} .

This invention thus proposes a new method for keyword spotting, based on recent advances in confidence measures, using local posterior probabilities, but without requiring the explicit use of filler models. A new method, referred to as Iterating Viterbi Decoding (IVD), to solve the above optimization problem with a simple DP process (not requiring to store pointers and scores for all possible ending and start times). Other three new beam-search [algorithms] *versions* corresponding to three different confidence measures are also proposed.

To summarize, the object of the invention consists of:

- Method of recognition of a subsequence using a direct maximization of confidence measures.
- The method of IVD for directly maximizing the confidence measures based on simple normalization.
- The use of the confidence measure and method of recognition named Real Fitting, based on individual fitting for each phoneme.
- Methods of recognition using simple and double normalization by:
 - combining these measures with additional confidence measures mentioned here, respectively the maximal length and real matching limitation.
- The use of the aforementioned methods in keyword recognition.
- The use of the aforementioned methods in subsequence recognition of organic matter.

- The use of the aforementioned methods in recognition of objects in images.

3 DETAILED DESCRIPTION OF PREFERRED EMBODIMENTS

Execution: The method can be performed using a personal computer or can be implemented in specialized hardware.

1. A representation under the form of an HMM is obtained for the subsequences that are looked for (word, protein profile, section of an image of the object).
2. A tool will be obtained (eventually trained Ex: for speech recognition) for the estimation of the posteriors. For example multi-Gaussians, neuronal networks, clusters, database with Generalized Profiles and mutation matrices (PAM, BLOSSUM, etc.).
3. One of the proposed [algorithms] *_versions_* should be implemented. They yield close performance but the method of Real Fitting coupled with a well checked dictionary should perform best.

For the first [algorithm] *_method_* (IVD)

- (a) The classic [algorithm of] = Viterbi is implemented with the modification that, for each pair $P = \langle sample, state \rangle$ one propagates the time-frame of transition between the state q_G and the states of the HMM M for the path that arrives at P. These are inherited from the path that wins the entrance in the pair P, excepting for the moment when their decision is taken, namely when they receive the index of the corresponding sample.

- (b) $w = -\log P(M|X_b^e)$ is computed by subtracting from the cumulated posterior that is returned by the Viterbi algorithm for the path $Q_{b_t}^{e_t}$, the value $(N - (e_t - b_t + 1)) * \varepsilon_t$ corresponding to the contribution of the states q_G and dividing the result through $e_t - b_t + 1$. $e_t - b_t + 1$ from the previous formula can be factored outside the fraction.
- (c) The initialization of ε is made with an expected mean value. One can use the w that is computed when the state q_G is associated with an emission posterior equal to the average of the best K emission probabilities of the current sample as done in the well-known garbage on-line model. In this case, K is trained using the corresponding technique.

The next Beam search [algorithms] *methods*, are implemented according to the description in the corresponding sections. For each pair $P = \langle sample, state \rangle$ one computes for each corresponding path the sum and length in the last phoneme, as well as the sum over the normalized cumulated posteriors of the previous phonemes (and their number). Also, the entrance and exit samples into the HMM M are computed and propagated like in the previous method, in order to ensure the localization of the subsequence.

4. If one searched entity (keyword, sequence, object) can have several HMM models, all of them are taken into consideration as competitors. This is the case of the words with several pronunciations (or of the objects that have different structures in different states, for the recognition in images).

After the computation of the confidence measure for each model of the subsequences,

one eliminates those with a confidence measure in disagreement with a threshold that is trained for the configuration and the goal of the given application. For example, for
3 speech recognition with neuronal networks and minus of the logarithm of the posteriors, the threshold is chosen in the wanted point of the ROC curve obtained in tests.

5. The remained alternatives are extracted in the order of their confidence measure and
6 with the elimination of the conflicting alternatives until exhaustion. Each time when an alternative is eliminated, the searched entity with the corresponding HMM is re-estimated for the remaining sections in the sequence in which the search is performed.
9 If the new confidence measure passes the test of the threshold, then it will be inserted in the position corresponding to its score in the queue of alternatives.

6. The successful alternatives can undergo tests of superior levels like for example a
12 question of confirmation for speech recognition, opinion of one operator, etc.

7. For objects recognition in images:

Posteriors are obtained by computing a distance between the color of the model and
15 that of element in the section of the image. If the context requires, the image will be preprocessed to ensure a certain normalization (Ex: changeable conditions of light will make necessary a transformation based on the histogram).

18 The phonemes of the speech recognition correspond to parts of the object. The structure (existence of transitions and their probabilities) can be modified, function of the characteristics detected along the current path. For example, after detecting regions
21 of the object with certain lengths, one can estimate the expected length of the remain-

ing regions. Thus, the number of the expected samples for the future states can be established and the HMM attached to the object will be configured accordingly.

- 3 A direction is scanned for the detection of the best fitting and afterwards, other directions will be scanned for discovering new fittings, as well as for testing the previous ones. The final test will be certified by classical methods such as cross-correlation or
- 6 by the analysis of the contours in the hypothesized position.

To mention some examples for the application of the proposed method:

- The recognition of keywords begins to be used in answering automates of banking
9 system as well as telephone and automates for control, sales or information. The method offers a possibility to recognize keywords in spontaneous speech with multiple speakers.
- 12 • The recognition of DNA sequences is important for the study of the human Genome. One of the biggest problem of the involved techniques consists in the high quantity of data that have to be processed.
- 15 • The recognition of objects in images is used, among others, in cartography and in the coordination of industrial robots. The method allows a quick estimation of the position of the objects in scenes and can be validated with extra tests, using classical methods
18 of cross-correlation.